

基于深度学习的面向 IP-over-EON 的可编程 跨层网络业务性能感知系统

朱祖勍, 孔嘉伟, 牛彬, 唐绍飞, 房红强, 刘思祺

(中国科学技术大学信息科学技术学院, 安徽 合肥 230027)

摘要: 为了实现实时的、细粒度的网络性能监测与调整, 并且满足不同应用的特定服务质量需求, 提出了基于深度学习的面向 IP-over-EON 的可编程跨层网络业务性能感知系统。该系统将基于网络业务性能感知的分布式网络监测与集中式网络管控相结合, 分布式网络监测实现跨层和细粒度的网络监控, 并基于深度学习进行数据分析。实验结果表明, 该系统通过有机地结合集中式与分布式的处理方式, 实现了及时的、自动化的网络控制与管理, 具有良好的可扩展性。

关键词: 深度学习; 弹性光网络; 跨层带内网络遥测; 网络异常监测与定位

中图分类号: TN919

文献标识码: A

doi: 10.11959/j.issn.1000-436x.2019229

DL-assisted programmable multilayer network application awareness system for IP-over-EON

ZHU Zuqing, KONG Jiawei, NIU Bin, TANG Shaofei, FANG Hongqiang, LIU Siqi

School of Information Science and Technology, University of Science and Technology of China, Hefei 230027, China

Abstract: In order to realize real-time and fine-granularity network monitoring and adjustment, and satisfy the specific QoS demands of various applications, a deep learning (DL) assisted programmable multilayer network application performance awareness system for IP-over-EON was proposed. The distributed network monitoring based on network application performance awareness was combined with centralized network management. The multilayer and fine-grained network monitoring was implemented by distributed network monitoring, and the data analysis through DL was performed. Experimental results indicate that by combining distributed and centralized processing seamlessly, the proposed network monitoring system can not only realize timely and automatic network control and management but also provide superior scalability.

Key words: deep learning, elastic optical network, multi-layer in-band network telemetry, network monitoring and failure localization

1 引言

当前, 随着云计算和大数据等新兴业务的飞速发展^[1-3], 主干网中的网络流量在飞速增长的同时, 其统计特性也变得越来越复杂。这些网络流量给现有基于 IP-over-Optical 的主干网体系架构带来了新的挑战^[4-6], 也使传统波分复用 (DWDM, dense

wavelength division multiplexing) 网络的可扩展性和灵活性呈现出不足。为了解决这一问题, 基于灵活栅格的弹性光网络 (EON, elastic optical network) 应运而生^[7-11], 它们提供更细的频谱分配粒度和更加灵活的频谱分配机制, 从而极大地提升了光层的灵活性和频谱利用率。在此基础上, 人们开始考虑用 EON 替换传统 IP-over-Optical 主干网中的

收稿日期: 2019-07-29; 修回日期: 2019-10-17

基金项目: 国家自然科学基金资助项目 (No.61871357, No.61771445, No.61701472)

Foundation Item: The National Natural Science Foundation of China (No.61871357, No.61771445, No.61701472)

DWDM 层, 以实现 IP-over-EON^[12-14], 集成 IP 技术和 EON 技术的优势。与此同时, 软件定义网络 (SDN, software-defined networking)^[15-16] 技术通过分离网络的控制面和数据面并提供集中式的网络控制与管理 (NC&M, network control and management) 机制来保障网络的可编程性, SDN 在 EON 中已经获得了应用并证明了其优越性^[17-20]。因此可以预见, 将 SDN 应用于 IP-over-EON 中, 即实现 SD-IPoEON (software-defined IP-over-EON), 有助于改善主干网的适配性和可编程性。

为了全面地体现 SD-IPoEON 架构的优势, 还需要实时的、细粒度的网络监测, 以保障其中各类网络应用的服务质量 (QoS, quality of service)。但是, 现有的传统网络监测技术难以满足实时监控和自动网络调整的需求, 造成这一问题原因主要有以下 3 个方面: 首先, 因为 IP-over-EON 是一个复杂的多层网络, 针对其的故障监测与定位需要采用跨层的方法, 同时分析 IP 层与 EON 层的情况, 然而, 当前的主流技术基本上都是单层的, 仅针对 IP 层或 EON 层进行分析, 并没有将两者结合起来; 其次, 现有技术大多仅通过带外的方式收集网络状态信息^[21-22], 这决定了它们难以实现实时的、端到端的和细粒度的监测, 而且它们会增大控制器的负担, 可扩展性较差; 最后, 现有技术较少考虑业务级别 (App-level, application level) 的监测, 考虑到不同的 App (Application), 有着各自不同的 QoS 需求, 它们的异常监测与故障恢复机制会有不同。

考虑到以上问题, 本文通过引入带内网络遥测 (INT, in-band network telemetry), 并在其基础上设计网络监控系统, 将 EON 层的信息通过 INT 的方式采集出来, 以实现跨层带内网络遥测 (ML-INT, multi-layer INT) 机制^[23], 同时将 IP 层和 EON 层的信息封装到 ML-INT 分组中, 完成实时的、分布式的和流粒度的端到端 App-level 网络监测。注意到, 尽管 App-level 的网络监控能很好地采集业务的端到端 QoS 参数, 它并不能单独完成网络故障的定位和网络状态的及时调整, 因此, 本文在这一分布式网络监控系统内加入集中式的网络监测, 利用 SDN 控制器完成链路级别 (link-level) 的网络监测。具体来说, App-level 监测通过 ML-INT 实现端到端的实时细粒度网络监测, 对采集到的数据进行分布式分析, 利用深度学习 (DL, deep learning) 实现故障的初步识别与定位^[24], 然后将异常信息上传给控制

器。控制器再结合自身 link-level 的带外监测信息, 进一步实现精确的故障定位并采取相应的故障恢复措施。本文提出的基于深度学习的跨层网络业务性能感知系统, 融合了分布式网络监测与集中式网络管控的优点。类似于人类的神经系统, 集中式的 SDN 控制器为大脑, 负责集中式网络管控; 分布式的 App-level 监测为神经感应元, 负责分布式的细粒度网络监测。因此, 将其称为网络神经系统 (NNS, network nervous system)^[25]。通过搭建小规模但真实的网络平台, 实现并用实验展示了 NNS 的优势, 实验结果证明其可以针对网络业务的 QoS 需求, 实现精确有效的故障识别、定位和恢复。同时, 实验结果也表明, 本文所提的 NNS 比传统的网络监测系统具有更良好的可扩展性和更灵活的可编程性。

2 系统架构

图 1 展示了 NNS 的整体架构^[23,25], 系统的数据平面由一个 IP-over-EON 组成。EON 层由光纤链路和带宽可变的光交换机 (BV-WSS, bandwidth variable wavelength-selective switch) 构成^[26-27], 并可以建立光路, 其中, 光性能监测器 (OPM, optical performance monitor) 监控各个链路的状态, 采集各条光路的光信噪比 (OSNR, optical signal-to-noise ratio)、中心波长和功率等信息。IP 层由可编程交换机和业务主机构成, 业务主机上运行着不同 QoS 需求的业务; 可编程交换机通过 INT 代理 (INT agent) 和 OPM 进行信息交互, INT agent 通过轮询的方式从 OPM 中获取所需的光路信息, 并发送至可编程交换机。如图 1 所示, 可编程交换机同时将 IP 层和 EON 层的数据插入 INT 头部的 INT 数据字段中。IP 层数据包括可编程交换机的交换机标识 (switch ID)、分组所经过的端口信息和排队时延信息, EON 层的数据包括 OSNR、中心波长和功率信息。当业务主机 A 与业务主机 B 的通信分组被转发到最后一台可编程交换机 (INT sink) 上, INT sink 将 INT 头部拆除并将其发送至数据分析设备 (data analyzer)。因此, ML-INT 技术对于通信的业务主机之间是透明的。data analyzer 通过解析 INT 头部的 INT 字段, 便可以依次获取业务主机 A 与 B 的通信路径上每一台可编程交换机和每一条光路的信息, 然后将解析出的 ML-INT 数据发送至 DL 故障识别和定位模块。

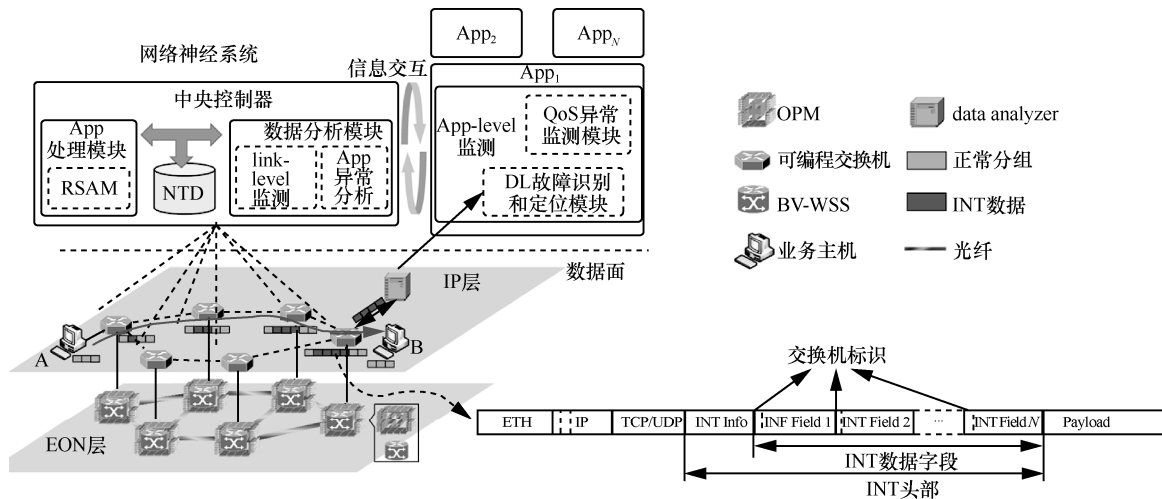


图 1 NNS 整体架构

由图 1 可知，App 中集成了 App-level 监测模块，App-level 监测模块包含 QoS 异常监测模块和 DL 故障识别和定位模块。其中，QoS 异常监测模块通过监测 QoS 参数，如接收端带宽和抖动等，可以实现基于业务的独特 QoS 需求定义异常，并及时将异常汇报给 DL 故障识别和定位模块。DL 故障识别和定位模块接收到异常信息后，开始对 data analyzer 发来的 ML-INT 数据进行分析，实现光路级别的故障识别和定位，并将初步的识别和定位结果传输至 SDN 中央控制器。中央控制器共包含 3 个模块，分别为数据分析模块、网络拓扑数据库(NTD, network topology database) 模块和 App 处理模块。数据分析模块包含 App 异常分析模块和 link-level 监测模块。App 异常分析模块负责接收 DL 故障识别和定位模块发送的故障初步识别和定位信息，link-level 监测模块通过带外遥测方式粗粒度地收集数据面信息。需要注意的是，DL 故障识别和定位模块是 App-level 的故障分析与定位，即光路级别的初步故障识别与定位。数据分析模块通过结合控制器收集到的 link-level 的光层信息和故障的初步识别和定位结果，实现 link-level 的精确定位。NTD 模块负责存储并更新数据平面网络拓扑信息。App 处理模块主要对业务的优先级和所属类别（例如是否时延敏感等）进行分析。当多个业务同时出现故障时，App 处理模块根据业务优先级设定故障恢复的先后顺序，并根据业务所属的类别，通过路由和频谱分配模块 (RSAM, routing and spectrum assignment module) 进行重路由（非时延敏感业务）或者切换到备份光路（时延敏感业务），以实现最

佳的故障恢复方案。

图 2 详细地说明了 NNS 在正常阶段和异常阶段时模块间的信息交互过程，其中模块 1 和模块 2 分别代表图 1 中 App-level 监测模块中的 QoS 异常监测模块和 DL 故障识别和定位模块。NNS 在正常阶段时，中央控制器通过带外遥测方式粗粒度的方式，即轮询时间间隔长，进行 link-level 监测。无论 NNS 处于正常阶段还是异常阶段，data analyzer 都会及时地将解析出的 ML-INT 数据发送至 DL 故障识别和定位模块。值得注意的是，只有在异常阶段下，当 DL 故障识别和定位模块接收到 QoS 异常监测模块的异常报告时，才会对 ML-INT 数据进行 App-level 的初步识别和定位，并将识别和定位结果报告至中央控制器的 App 异常分析模块。中央控制器收到异常报告后，将会从粗粒度方式变化为细粒度方式，即轮询间隔时间短，有针对性地对定位的故障光路进行 link-level 的监测。NNS 通过结合 DL 故障识别和分类模块 App-level 的初步识别和定位结果和控制器 link-level 的监测信息进行精确的故障定位。之后中央控制器通过 App 处理模块获得业务的优先级和相关信息，根据 NTD 模块中存储的数据面拓扑信息进行 IP 层重路由或光路重配置等 App-level 的故障恢复方案。

3 系统实现

3.1 可编程 ML-INT

为了实现 ML-INT 技术，通过可编程协议无关分组处理 (P4, programming protocol-independent packet processor) 语言^[28]在可编程交换机上对数据

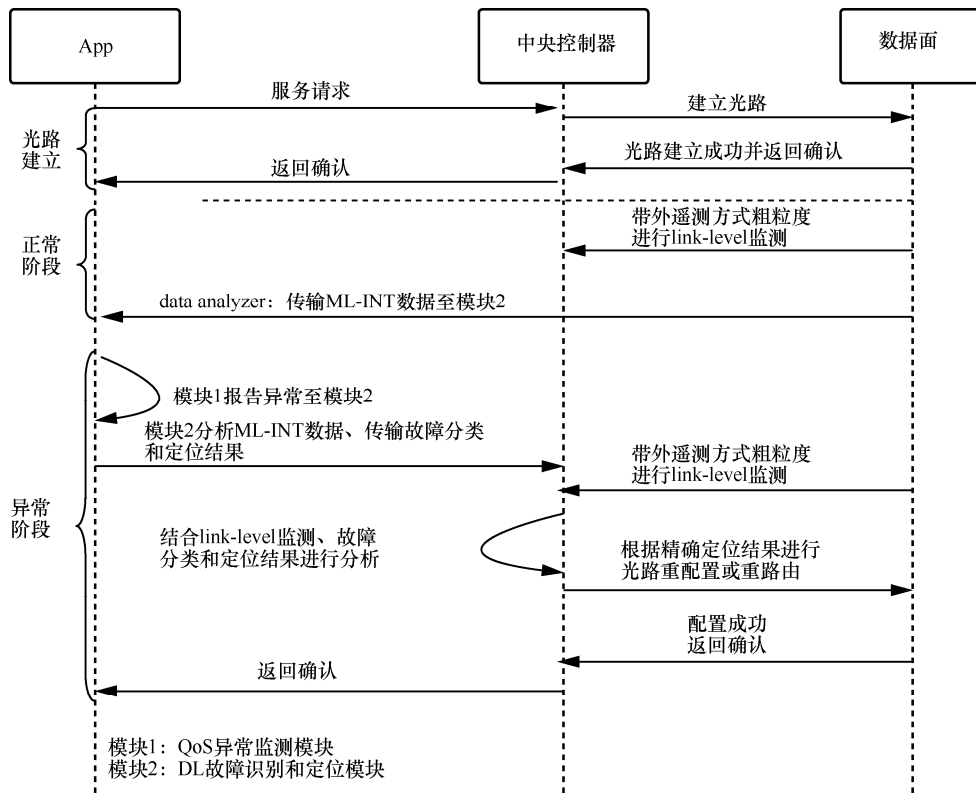


图 2 NNS 交互流程

包的格式进行相关的处理。首先，当业务分组经过第一台可编程交换机（INT source）时，INT source 以一定的比例（例如 50%）在业务分组中封装 INT 头部，同时定义 INT 结构并在 INT 头部定义 INT 数据字段。IP 层和 EON 层数据通过可编程交换机写入 ML-INT 分组中的 INT 数据字段。接着，ML-INT 分组被转发至下一台可编程交换机（INT intermediate hop），INT intermediate hop 将继续向 ML-INT 分组写入本机相关的 IP 层信息和 EON 信息。最后，INT sink 将 INT 头部拆除，然后发送至 data analyzer 进行解析。为了尽量减小 ML-INT 所带来的开销，本文指定每个可编程交换机在经过的 ML-INT 分组中只插入 2 个 INT 数据字段，从而减小 INT 头部的长度。2 个 INT 数据字段中，第一个 INT 数据字段固定为 switch ID，指示数据分组的转发路径和所经过的交换机信息；第二个 INT 数据字段为 IP 层排队时延、EON 层 OSNR 或功率，可编程交换机将这些数据依次写入多个 ML-INT 分组中。由于网络中业务通信速率通常较高（例如 10 Gbit/s），分组之间的时间间隔都是微秒级别。通常情况下，网络状态难以在微秒级别的时间内发生剧变，因此将

IP 层和 EON 层的数据分多次写入 ML-INT 分组中是合理的。与此同时，data analyzer 需要解析多个 ML-INT 分组方可获取完整的端到端 ML-INT 数据。

3.2 基于深度学习的 App-level 监测

NNS 架构中，DL 故障识别和定位模块主要负责对 data analyzer 解析出的 ML-INT 数据进行初步故障识别和定位。传统的网络架构通常基于阈值的方式判断某类故障是否发生。随着网络规模的不断发展，故障类型逐渐增多，为不同的故障设置不同的阈值将会增加网络操作复杂度，可扩展性较差。其次，有的故障类型难以仅通过简单的阈值进行识别。如图 3 所示，本文通过在光路中引入噪声产生故障，当接收端带宽低于发送端的一定比例（例如 90%）时即判定为异常。从图 3 中可以看出，第 3 s 的功率值低于第 5 s 的功率值，然而第 3 s 处于正常阶段，而第 5 s 处于异常阶段。这是因为引入噪声对功率影响较小，但当 OSNR 降低到一定值时，也会引发异常。所以对于噪声故障需要综合考虑功率和 OSNR，不能只通过阈值进行判断。因此，本文采取基于深度学习的 DL 故障识别和分类方法。

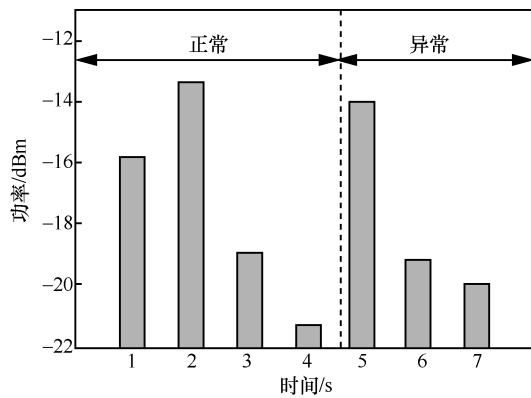


图 3 App-level 功率监测

DL 故障识别和分类模块是深度神经网络 (DNN, deep neural network) 架构, 包含输入层、两层隐藏层和输出层, 通过监督式学习进行训练。输入层节点数等于要处理的数据的变量数, 即 3 个; 两层隐藏层的节点数都设置为 128; 输出层节点数等于每个输入对应的输出数, 为 4。DNN 架构中输出层的激活函数为 softmax 函数, 并选取分类问题中常用的交叉熵函数作为损失函数来刻画预测的概率分布和真实输出的概率分布之间的距离。DNN 训练时通过反向传播和梯度下降算法调整神经网络中参数的取值, 从而最小化损失函数。DL 故障识别和分类模块的输入是 ML-INT 数据, 包括 IP 层排队时延、EON 层的功率和 OSNR。值得注意的是, 若 ML-INT 分组经过 N 个可编程交换机, 端到端完整 ML-INT 数据会包含 N 组输入 (一组输入为排队时延、功率和 OSNR)。data analyzer 在将端到端完整 ML-INT 数据发送至 DL 故障识别和定位模块时, 会先依次从业务分组的发送端到目的端对 ML-INT 数据进行 N 组划分。DL 故障识别和定位模块依次对 N 组数据进行故障识别, 第一个出现故障的光路或可编程交换机即为故障位置, 从而实现故障定位, 这种每次仅对一组数据进行分析的方式适用于不同的网络规模, 具有很好的可扩展性。

本文通过搭建一个小规模但是真实的 IP-over-EON 测试平台, 模拟产生各种故障 (拥塞异常、功率衰减异常和噪声异常), 从而获得训练数据。例如, 本文通过功率衰减器衰减光路中的功率, 在光路中引入噪声和加入背景流让可编程交换机产生拥塞来分别模拟功率衰减异常、噪声异常和拥塞异常。同时, 如果接收端的带宽低于发送端的一定比例 (例如 90%), 即判定产生异常, 为训练数据贴上标签 (例如 0、1、2 和 3 分别代表正常、功率衰

减异常、噪声异常和 IP 层拥塞异常)。训练时, 每次只产生一种异常。故障样本总共采集了约 17 000 组, 每组训练数据包含排队时延、功率和 OSNR 数据, 其中训练集和测试集各占 90%和 10%。通过离线训练, AI 模型在测试集上的准确率达到 96.99%。

4 实验展示与结果分析

4.1 NNS 实验验证

本文通过搭建如图 4 所示的实验平台展示和验证 NNS 架构功能。IP 层包含 6 台可编程交换机, 端口为 10 GbE 光端口, 通过 2 台商用的分组发送工具来模拟端上的 2 个业务主机进行通信。EON 层包括 6 个 1×9 的 BV-WSS 光纤、6 个 OPM 和 8 个掺铒光纤放大器 (EDFA, erbium-doped fiber amplifier)。BV-WSS 工作频率范围为 1 528.43~1 566.88 nm, 频谱分配粒度为 12.5 GHz。OPM 采用光信息监测仪 (OCM, optical channel monitor), 光谱分辨率为 312.5 MHz。图 4 中的箭头表示的是业务通信时的路由路径。从图 4 中可以看出, 业务主机 A 与业务主机 B 通信时经过 2 个光路。第一个光路从可编程交换机 1 到可编程交换机 2, 包含 a 一条链路, 第二个光路从可编程交换机 2 到可编程交换机 3, 包含 b 和 c 共两条链路。中央控制器基于开放式网络操作系统 (ONOS, open network operating system) 平台与数据平面进行通信。为了验证 NNS 架构的有效性, 本文在 EON 层和 IP 层各做了一个故障识别、定位和恢复的实验。实验中发送端业务的吞吐量为 8 Gbit/s, 分组大小为 1 024 B, INT source 上 ML-INT 分组的比例设置为 50%。EON 层的实验中, 通过功率衰减器在 c 链路上对功率进行衰减, 用 Wireshark 抓取 NNS 架构间的通信数据分组, 从而验证 NNS 架构的功能。图 5(a)虚线框展示的是 data analyzer (IP 地址为 192.168.108.40) 向 DL 故障识别和定位模块 (IP 地址 192.168.108.229) 发送的一组 ML-INT 数据, 包括 switch ID、排队时延、功率和 OSNR。每个可编程交换机都有预先定义好的 switch ID, 通过 switch ID 对故障进行定位。例如, 业务主机 A 和业务主机 B 通信时 IP 层经过 3 个可编程交换机, switch ID 分别为 1、2、3。在 c 链路进行功率衰减, 第一条光路正常, 第二条光路异常。DL 故障识别和定位模块依次 ML-INT 数据进行分析, 故障类型为 1, switch ID 在 2 和 3 之间。如图 5(b)虚线框所示,

Wireshark 抓取的是 DL 故障识别和定位模块发送至中央控制器 (IP 地址为 192.168.108.225) 的数据分组, 分组中包含的故障类型为 1, switch ID 为 2 和 3。中央控制器获取异常信息后, 通过带外遥测方式细粒度 (例如每 1 s 获取 1 次) 获取 switch ID 在 2 和 3 之间所有链路的光路信息 (中心频率、功率和 OSNR) 并与正常阶段下通过粗粒度 (例如每 10 s 获取 1 次) 获取的光路数据进行对比和分析。由于 DL 故障识别和定位模块传输故障类别为 1, 即功率衰减故障, 中央控制器依次分析 b 链路和 c 链路最近一段时间的功率变化趋势, 通过功率变化范围判断 (例如功率降低范围超过 4 dBm) 链路是否产生异常, 从而得出 c 链路故障。最后, 中央控制器通过 App 处理模块获取业务的优先级和所属服务类型等信息, 将 c 链路通过 BV-WSS 切换到备份链路进行故障恢复。

接收端统计的是每秒收到的帧数, 当数据分组大小为 1 024 B 时, 8 Gbit/s 约为 950 000 frame/s。当 c 链路功率衰减异常时, 接收端的带宽降为约 600 000 frame/s, 低于 QoS 异常监测模块的阈值。c 链路上 OPM 监测的功率随时间变化趋势如图 6(a) 所示。OPM 上功率恢复时延经过 5 次测量取平均值, 约为 1.97 s。恢复时延主要包括 OPM 扫描光谱、DL 故障识别和定位模块进行故障初步识别和定位、中央控制器下发光流表和 BV-WSS 重新配置切换链路的时延。图 6(b) 展示的是业务主机 B 的接受带宽随时间变化趋势, 端到端恢复时延约为 2.99 s, 之所以比 OPM 上功率的恢复时延要长, 主要是由于本文通过商用的分组发送工具来模拟端上的 2 个业务主机进行通信, 然而分组发送工具接收端带宽以秒级进行更新, 会对端到端的恢复时延造成一定的影响。

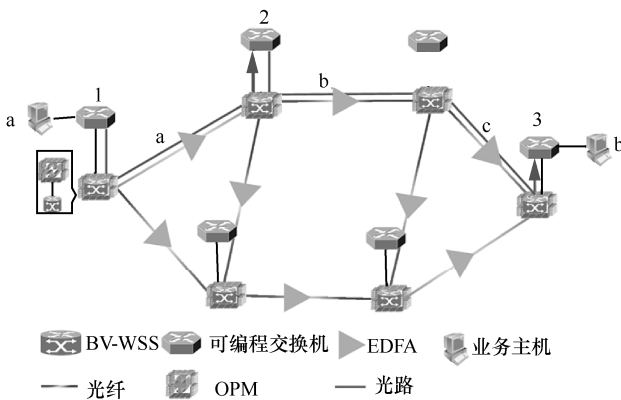


图 4 实验测试平台

在 IP 层, 通过加背景流的方式使 Switch ID 为 2 的可编程交换机产生拥塞。背景流发送速率为 4 Gbit/s, 业务主机 A 和业务主机 B 之间通信速率为 8 Gbit/s, 可编程交换机的端口最大接受速率为 10 Gbit/s。因为 IP 层的拥塞会导致分组丢失, 从而导致接收端带宽变低, 所以 QoS 异常监测模块通过监测接收端的带宽来判断是否有异常发生, 然后将异常信息报告至 DL 故障识别和分类模块。DL 故障识别和分类模块接收到异常信息后, 对 ML-INT 数据进行分析, 并将故障识别和定位结果报告至中央控制器, 其数据分组和图 5(b) 所示的数据格式相同,

14500	95.939992235	192.168.108.40	192.168.108.229	TCP	82 ML-INT Data
14600	95.939259042	192.168.108.229	192.168.108.40	TCP	66 34296 → 8887 [ACK] Seq=1 Ack=1327425 Win=184832 Len=0 TSval=2567323157 TSecr=3254287184
14601	95.939767528	192.168.108.40	192.168.108.229	TCP	82 ML-INT Data
14602	95.939803052	192.168.108.229	192.168.108.40	TCP	66 34296 → 8887 [ACK] Seq=1 Ack=1327441 Win=184832 Len=0 TSval=2567323161 TSecr=3254287185
14603	95.939995235	192.168.108.40	192.168.108.229	TCP	82 ML-INT Data

```

Frame 14599: 82 bytes on wire (656 bits), 82 bytes captured (656 bits) on interface 0
Ethernet II, Src: WstronI_0d:a:e:8c (70:e2:84:0d:a:e:8c), Dst: SuperMic_81:66:6e (ac:1f:6b:81:66:6e)
Internet Protocol Version 4, Src: 192.168.108.40 [192.168.108.40], Dst: 192.168.108.229 [192.168.108.229]
Transmission Control Protocol, Src Port: 8887 (8887), Dst Port: 34296 (34296), Seq: 1327409, Ack: 1, Len: 16
ML-INT Data
  SwitchID: 0xf1000000
  Delays: 0x60000000
  Power: 0x8fffffff
  OSNR: 0x15010000
  
```

(a) data analyzer 传输 ML-INT 数据至 DL 故障识别和定位模块

2353	29.416921937	192.168.108.229	192.168.108.225	TCP	110 47662 → 9038 [ACK] Seq=1 Ack=45 Win=29312 Len=44 TSval=3068399908 TSecr=260660132
2354	29.416954975	192.168.108.225	192.168.108.229	TCP	66 9038 → 47662 [ACK] Seq=1 Ack=45 Win=29056 Len=0 TSval=2806791111 TSecr=3068399908
2355	29.443371071	192.168.108.229	192.168.108.25	TCP	54 36298 → ddi-tcp-1(8888) [ACK] Seq=1 Ack=133 Win=29312 Len=0

```

Frame 2353: 110 bytes on wire (880 bits), 110 bytes captured (880 bits) on interface 0
Ethernet II, Src: SuperMic_81:66:6e (ac:1f:6b:81:66:6e), Dst: Inventec_e5:74:70 (7c:d3:0a:e5:74:70)
Internet Protocol Version 4, Src: 192.168.108.229 [192.168.108.229], Dst: 192.168.108.225 [192.168.108.225]
Transmission Control Protocol, Src Port: 47662 (47662), Dst Port: 9038 (9038), Seq: 1, Ack: 1, Len: 44
Data (44 bytes)
  Data: 7b 22661756c74223a20312c2022737697463869643122...
  (Length: 44)
  
```

```

1000 7c d3 0a e5 74 70 ac 1f 6b 81 66 6e 08 00 45 00 |...tp..k.fn..E.
1010 00 60 dd c9 40 00 40 06 01 b7 c0 a8 6c e5 c0 a8 |...@. ....l...
1020 6c e1 ba 2e 23 4e 0e d4 a3 4c 21 9d de bd 80 18 |...#N...!l.....
1030 00 e5 5b 6a 00 00 01 01 08 0a b8 b5 5e 54 10 ba |...[...-I...
1040 a0 14 7b 22 66 61 75 6c 74 22 3a 20 31 2c 20 22 |..(*'faul t': 1, ")
1050 73 77 69 74 63 68 69 64 31 22 3a 20 32 2c 20 22 |switchid 1': 2, "1
1060 73 77 69 74 63 68 69 64 32 22 3a 20 33 7d |switchid 2': 3)
  
```

(b) DL 故障识别和定位模块传输初步识别和分类结果至中央控制器

图 5 Wireshark 抓取数据分组验证 NNS 架构

故障类型为电层拥塞（对应数字为 3），switch ID 1 和 2 传输的分别为拥塞交换机的上一跳可编程交换机（switch ID 为 1）和拥塞的可编程交换机（Switch ID 为 2）的 switch ID。中央控制器结合 NTD 存储的底层拓扑信息，通过 RSAM 从上一跳可编程交换机进行重路由，避免拥塞的可编程交换机。

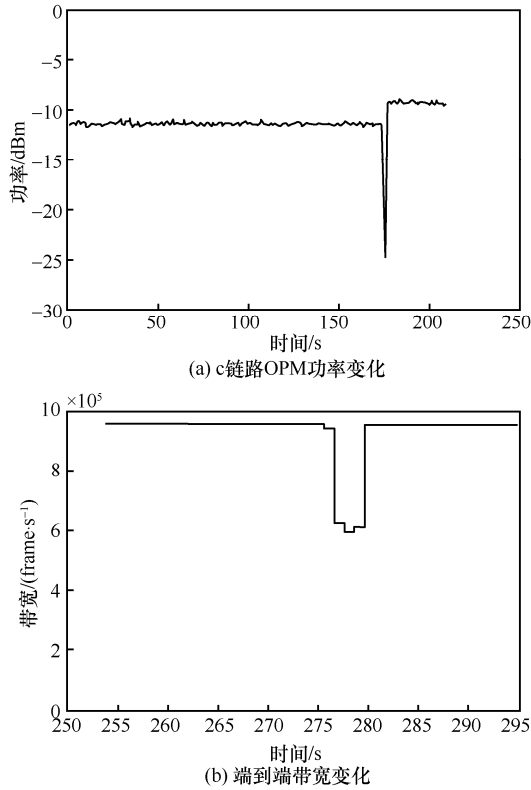


图 6 EON 层故障分类、定位和恢复

图 7(a)展示的是 switch ID 为 2 的交换机上排队时延变化趋势，正常阶段下排队时延约为 0.1 μs ，拥塞时最高排队时延达到 120 μs 以上。图 7(b)展示的是接收端带宽变化趋势。经过 5 次测量求平均值，可编程交换机上的排队时延恢复时间约为 0.85 s，接收端的带宽恢复时延约为 1.99 s。由于分组发送工具接收端的带宽数据以秒级进行更新，导致端到端的带宽恢复时延长于可编程交换机的恢复时延。

4.2 NNS 架构可扩展性验证

为了进一步展示 NNS 架构带来的可扩展性优势，如图 8 所示，通过仿真的方式产生多组数据，横坐标是数据量，每组数据包括时延、功率和 OSNR；纵坐标表示对应数据量下所需要的故障分类时间。分类时间在显卡为 GeForce GTX 1080Ti 的 GPU 上进行测试。当数据量低于 100 组时，分类时间约为 0.2 s，当数据量高于 10 000 组时，分类

时间明显增多。NNS 采取的是基于 ML-INT 的分布式网络数据分析，神经网络所需的数据量基本在 100 组以下。传统的纯带外遥测集中式地对数据面数据进行分析，当网络设备和链路增多，或者需要采取比链路级别更加精细的网络监测与定位时（例如需要精确定位到链路中具体的设备是否故障，需要监测的数据量进一步增加），都会显著增加神经网络的分类时间和定位时间。值得注意的是，因为纯带外收集数据的方式没有结合 App-level 端到端的异常信息，需要时刻对获取到的数据进行分析；而 NNS 架构在无异常发生时，不需要对收集到的数据进行分析，仅需在产生异常时对数据进行初步识别和定位，因此 NNS 架构在一定程度上减轻了神经网络的处理负担。

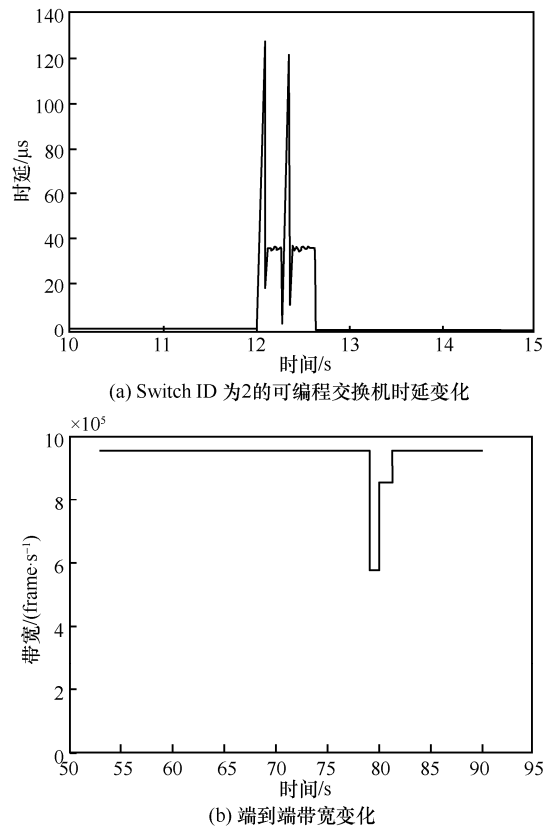


图 7 IP 层故障分类、定位和恢复

NNS 架构采取 App-level 的分布式方式分析，因此无异常发生时，中央控制器进行粗粒度的带外遥测，异常发生时再对故障光路进行细粒度的带外遥测。纯带外遥测方式的传统架构为了实现细粒度和实时的网络监测，无论是正常还是异常阶段，中央控制器都需要进行细粒度的带外遥测。图 9 展示了 NNS 架构和传统架构在正常情况和异常情况下

南向通信量的比例。本文设定粗粒度的轮询时间间隔为 10 s，细粒度的轮询时间间隔为 1 s。为了确保通信分组的大小相同，本文用同一个 OPM 进行实验，通过 Wireshark 获取中央控制器一段时间(10 min)内接收到的平均分组速度。如图 9 所示，通过 5 次实验获取平均值，正常阶段下，NNS 架构与传统架构的每秒平均分组速度比为 0.379:2.359；异常阶段下，两架构每秒的平均分组速度之比为 1:1。当发生异常的频率较低时，NNS 中央控制器每秒平均分组速度将远低于传统架构，进一步减少了控制器的负担，增加了可拓展性。

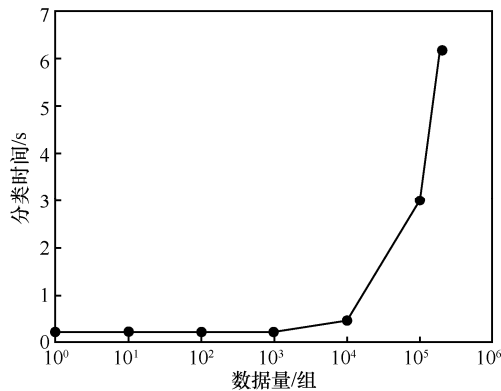


图 8 AI 模块分类时间趋势

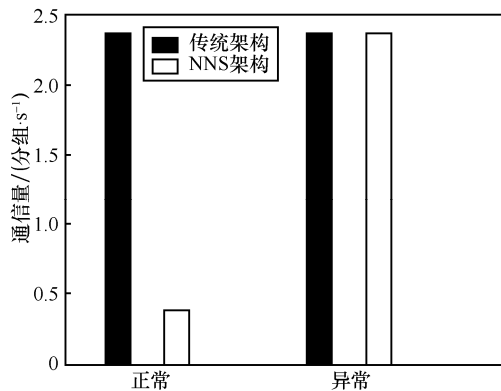


图 9 中央控制器南向通信量对比

5 结束语

本文提出了基于深度学习的跨层网络业务性能感知系统，融合了分布式网络监测与集中式网络管控的优点，实现了细粒度和实时的网络监测。本文通过搭建小规模实验平台，实现并展示了 NNS 架构的优势。无论 IP 层还是 EON 层产生异常，NNS 架构都可以精确快速地进行故障分类、定位和恢复。同时，通过故障识别和定位时间与

中央控制器南向通信量的比较，NNS 架构相比传统的纯带外监测系统，有着更加良好的可拓展性。

参考文献:

- [1] LU P, ZHANG L, LIU X H, et al. Highly-efficient data migration and backup for big data applications in elastic optical inter-datacenter networks[J]. IEEE Network, 2015, 29(5): 36-42.
- [2] KONG B X, LIU S Q, YIN J, et al. Demonstration of application-driven network slicing and orchestration in optical/packet domains: on-demand vDC expansion for Hadoop MapReduce optimization[J]. Optics Express, 2018, 26(11): 14066-14085.
- [3] LU P, SUN Q Y, WU K Y, et al. Distributed online hybrid cloud management for profit-driven multimedia cloud computing[J]. IEEE Transaction on Multimedia, 2015, 17(8): 1297-1308.
- [4] YAO J J, LU P, GONG L, et al. On fast and coordinated data backup in geo-distributed optical inter-datacenter networks[J]. Journal of Lightwave Technology, 2015, 33(14): 3005-3015.
- [5] LU P, ZHU Z Q. Data-oriented task scheduling in fixed- and flexible-grid multilayer inter-DC optical networks: a comparison study[J]. Journal of Lightwave Technology, 2017, 35(24): 5335-5346.
- [6] XIE X K, LING Q, LU P, et al. Evacuate before too late: distributed backup in inter-DC networks with progressive disasters[J]. IEEE Transaction on Parallel and Distributed Systems, 2018, 29(5): 1058-1074.
- [7] GERSTEL O, JINNO M, LORD A, et al. Elastic optical networking: a new dawn for the optical layer[J]. IEEE Communications Magazine, 2012, 50(2): s12-s20.
- [8] ZHU Z Q, LU W, ZHANG L, et al. Dynamic service provisioning in elastic optical networks with hybrid single-multi-path routing[J]. Journal of Lightwave Technology, 2013, 31(1): 15-22.
- [9] GONG L, ZHOU X, LIU X H, et al. Efficient resource allocation for all optical multicasting over spectrum-sliced elastic optical networks[J]. Journal of Optical Communication and Networking, 2013, 5(8): 836-847.
- [10] YIN Y W, ZHANG H, ZHANG M Y. Spectral and spatial 2D fragmentation aware routing and spectrum assignment algorithms in elastic optical networks[J]. Journal of Optical Communication and Networking, 2013, 5(10): A100-A106.
- [11] GONG L, ZHU Z Q. Virtual optical network embedding (VONE) over elastic optical networks[J]. Journal of Lightwave Technology, 2014, 32(3): 450-460.
- [12] GERSTEL O, FILSFILS C, TELKAMP T, et al. Multi-layer capacity planning for IP-optical networks[J]. IEEE Communication Magazine, 2014, 52(1): 44-51.
- [13] LIU S Q, LU W, ZHU Z Q. On the cross-layer orchestration to address IP router outages with cost-efficient multilayer restoration in IP-over-EONs[J]. Journal of Optical Communication and Networking, 2018, 10(1): A122-A132.
- [14] LU W, YIN X F, CHENG X B, et al. On cost-efficient integrated multilayer protection planning in IP-over-EONs[J]. Journal of Lightwave Technology, 2017, 36(10): 2037-2048.
- [15] FEAMSTER N, REXFORD J, ZEGURA E. The road to SDN: an intellectual history of programmable networks[J]. ACM SIGCOMM Computer Communication Review, 2014, 44(2): 87-98.

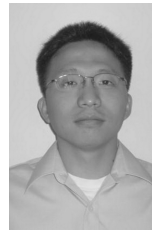
- [16] LI S R, HU D Y, FANG W J, et al. Protocol oblivious forwarding (POF): software-defined networking with enhanced programmability[J]. IEEE Network, 2017, 31(2): 58-66.
- [17] CHEN C, CHEN X L, ZHANG M Y, et al. Demonstrations of efficient online spectrum defragmentation in software-defined elastic optical networks[J]. Journal of Lightwave Technology, 2014, 32(24): 4701-4711.
- [18] ZHU Z Q, CHEN X L, CHEN C, et al. OpenFlow-assisted online defragmentation in single-multi-domain software-defined elastic optical networks[J]. Journal of Optical Communication and Networking, 2015, 7(1): A7-A15.
- [19] CHEN X L, TORNATORE M, ZHU S L, et al. Flexible availability-aware differentiated protection in software-defined elastic optical networks[J]. Journal of Lightwave Technology, 2015, 33(18): 3872-3882.
- [20] ZHU Z Q, CHEN C, CHEN X L, et al. Demonstration of cooperative resource allocation in an OpenFlow-controlled multidomain and multinational SD-EON testbed[J]. Journal of Lightwave Technology, 2015, 33(8): 1508-1514.
- [21] YAN S Y, AGUADO A, OU Y, et al. Multilayer network analytics with SDN-based monitoring framework[J]. Journal of Optical Communication and Networking, 2017, 9(2): A271-A279.
- [22] PAOLUCCI F, SGAMBELLURI A, CUGINI F, et al. Network telemetry streaming services in SDN-based disaggregated optical networks[J]. Journal of Lightwave Technology, 2018, 36(15): 3142-3149.
- [23] NIU B, KONG J W, TANG S F, et al. Visualize your IP-over-optical network in realtime: a P4-based flexible multilayer in-band network telemetry (ML-INT) system[J]. IEEE Access, 2019, 7: 82413-82423.
- [24] CHEN X L, LI B J, PROIETTI R, et al. Self-taught anomaly detection with hybrid unsupervised/supervised machine learning in optical networks[J]. Journal of Lightwave Technology, 2019, 37(7): 1742-1749.
- [25] FANG H Q, LU W, LIANG L P, et al. Building network nervous system with multilayer telemetry to realize AI-assisted reflexes in software-defined IP-over-EONs for application-aware service provisioning[C]//Optical Fiber Communication Conference. 2019: 1-3.
- [26] YIN J, GUO J N, KONG B X, et al. Experimental demonstration of

building and operating QoS-aware survivable vSD-EONs with transparent resiliency[J]. Optics Express, 2017, 25(13): 15468-15480.

- [27] ZHU Z Q, KONG B X, YIN J, et al. Build to tenants' requirements: On-demand application-driven vSD-EON slicing[J]. Journal of Optical Communication and Networking, 2018, 10(2): A206-A215.

- [28] BOSSHART P, DALY D, GIBB G, et al. P4: programming protocol independent packet processors[J]. ACM SIGCOMM Computer Communication Review, 2014, 44(3): 87-95.

[作者简介]



朱祖劼 (1979-)，男，安徽蚌埠人，博士，中国科学技术大学教授，主要研究方向为下一代网络体系架构、云计算和光通信网络以及多媒体网络。

孔嘉伟 (1995-)，男，安徽六安人，中国科学技术大学硕士生，主要研究方向为光通信网络。

牛彬 (1994-)，男，安徽阜阳人，中国科学技术大学硕士生，主要研究方向为下一代网络体系架构。

唐绍飞 (1995-)，男，安徽铜陵人，中国科学技术大学硕士生，主要研究方向为下一代网络体系架构。

房红强 (1995-)，男，安徽阜阳人，中国科学技术大学硕士生，主要研究方向为光通信网络和机器学习。

刘思琪 (1993-)，男，河南新乡人，中国科学技术大学博士生，主要研究方向为光通信网络和机器学习。

优青人物介绍

朱祖劼，教授，2014 年获得中国科学院“卓越青年科学家计划”支持，2018 年获中组部“青年拔尖人才计划支持”。发表 SCI 论文 110 余篇，曾获 ICC 2013、GLOBECOM 2013、ICNC 2014、ICC 2015 和 ONDM 2018 等旗舰学术会议的最佳论文奖；担任 IEEE CM、TNSM、OE、OSN 等多个期刊的 Series Editor 或 Editor。当前是 IEEE HPSR 会议的 Steering Committee Member，入选 IEEE Communications Society Distinguished Lecturer (2018—2019)。2011 年入选 IEEE 学会高级会员，2013 年入选美国光学学会 (OSA) 高级会员。